

STICHTING
MATHEMATISCH CENTRUM
2e BOERHAAVESTRAAT 49
AMSTERDAM

ZW 1950-010

Over "upcrosses" en "downcrosses" in tijdreeksen

"Actualiteiten"

H. Theil



1950

Voordracht door H. Theil in de serie
Actualiteiten op 20 Mei 1950.

Over "upcrosses" en "downcrosses" in tijdreeksen.

1. Onder tijdreeks verstaan we hier een getallenrij x_t , gedefinieerd op een discrete verzameling van waarden van de parameter t . Deze parameter wordt als de tijd geïnterpreteerd. We beperken ons tot het geval, dat de waarden van t aequidistant zijn; zonder verlies van algemeenheid kunnen we dan het verschil = 1 stellen. Verder beperken we ons tot waarden $t \geq t_0$, waarbij we (zolang t_0 constant gehouden wordt) zonder beperking kunnen aannemen, dat $t_0 = 0$ is.

In de literatuur worden ook algemene tijdreeksen beschouwd, bijv. die welke gedefinieerd zijn voor alle reële waarden van t .

2. Beschouw in het bijzonder de tijdreeksen, gegeven door een lineaire differentievergelijking met constante coëfficiënten. In deze van de orde h , dan heeft deze de vorm

$$(1) \quad x_t + a_1 x_{t-1} + \dots + a_h x_{t-h} = \varepsilon_t \quad (a_h \neq 0)$$

Hierbij is ε_t een gegeven tijdreeks en zijn a_1, \dots, a_h gegeven getallen. Meer in het bijzonder beschouwen we het geval $h = 1$

$$(2) \quad x_t = \alpha x_{t-1} + \varepsilon_t$$

Indien geldt: $\varepsilon_t \equiv 0$, dan hebben we een homogene differentievergelijking, waarvan de oplossing luidt

$$(3) \quad x_t = \alpha^t x_0.$$

Indien geldt: $|\alpha| < 1$, hebben we: $\lim_{t \rightarrow \infty} x_t = 0$.

De inhomogene differentievergelijking heeft tot oplossing:

$$(4) \quad x_t = \sum_{j=0}^t \alpha^{t-j} \varepsilon_j + \alpha^t x_0$$

indien we $x_0 = \varepsilon_0$ stellen, kunnen we voor (4) schrijven:

$$(5) \quad x_t = \sum_{j=0}^t \alpha^{t-j} \varepsilon_j.$$

3. Op allerlei gebieden van toegepaste wiskunde komen tijdreeksen voor bijv. op meteorologisch gebied (de dagelijkse of jaarlijkse regenval op een bepaalde plaats, enz), op agrarisch gebied (de gemiddelde jaarlijkse opbrengsten per HA van een bepaald gewas in een be

paald gebied, enz.), op demografisch gebied (de loop van een bevolking, geboortecijfers, enz.), op economisch gebied (maandelijke productie-indices, enz.). De theorie van de tijdreeksen stelt zich ten doel, deze empirische reeksen wiskundig te analyseren. Dit doel zou bereikt zijn, indien bijv. de empirische reeks een oplossing bleek te zijn van een differentievergelijking. Aangezien het aantal waarnemingen van dergelijke reeksen steeds eindig is, is het altijd mogelijk tot een dergelijke differentievergelijking te komen. Echter, wanneer men op grond hiervan voorspellingen doet, komen zij (behalve in de astronomie of de natuurkunde) doorgaans niet uit.

Hieraan kan tegemoet gekomen worden, doordat we de empirische reeks beschouwen als een "toevallige" keuze uit een verzameling van eventuele tijdreeksen. We nemen dus aan, dat voor iedere t \underline{x}_t een stochastische grootheid ¹⁾ is, d.w.z. dat \underline{x}_t een waarschijnlijkheidsverdeling heeft. \underline{x} is dus een stochastische functie van t . Op de algemene theorie van de stochastische functies en de stochastische processen gaan we hier niet in.

We zullen ons ook hier beperken tot het geval, dat een tijdreeks kan worden voorgesteld door een differentievergelijking van de eerste orde met constante α . $\underline{\varepsilon}_t$ is dan een stochastische functie van t . We nemen aan, dat $\underline{\varepsilon}_0, \dots, \underline{\varepsilon}_t$ onafhankelijk verdeeld zijn volgens dezelfde verdeling.

De oplossing (5) wordt nu

$$(6) \quad \underline{x}_t = \sum_{j=0}^t \alpha^{t-j} \underline{\varepsilon}_j.$$

4. Bij gegeven α en gegeven verdeling van $\underline{\varepsilon}_t$ kan men de verdeling vinden van \underline{x}_t of van bepaalde functies daarvan. Nu heeft men in empirische reeksen vaak gezocht naar periodiciteiten, bijv. met behulp van de zgn. periodigram-analyse. Voor een onderzoek naar deze periodiciteiten zullen wij ons hier echter richten naar Kendall's begrippen "upcrosses" en "downcrosses". We nemen daartoe aan, dat de verwachting van $\underline{\varepsilon}_t$ en dus ook van alle \underline{x}_t gelijk 0 is. Dan is een "upcross" resp. een "downcross" een punt $t+1$ langs de tijdas, waarvoor geldt

$$(7) \quad \begin{array}{ll} \underline{x}_t < 0 & \underline{x}_t > 0 \\ \text{resp.} & \\ \underline{x}_{t+1} > 0 & \underline{x}_{t+1} < 0 \end{array}$$

¹⁾ Stochastische grootheden worden door onderstreping onderscheiden van grootheden, die geen verdeling hebben.

We veronderstellen, dat de verdelingsfunctie van ε_t en dus van alle x_t continu is. Dan is $P[x_t = 0] = 0$. Hierdoor sluiten we uit de mogelijkheden: $x_{t-1} > 0 = x_t < x_{t+1}$ en $x_{t-1} < 0 = x_t > x_{t+1}$.

Nu volgt uit (6), dat het zinvol is, te spreken over de waarschijnlijkheid, dat voor een tweetal punten $t, t+1$ de ongelijkheden (7) gelden. Langs de tijdas zal men dan een verdeling van punten t , die "upcrosses" resp. "downcrosses" zijn, vinden. De afstand tussen twee opeenvolgende van dergelijke punten is dus een stochastische grootheid. We zullen ons nu bezig houden met de verdeling van deze afstand.

5. M.G. Kendall²⁾ heeft zich bezig gehouden met de verwachting van deze afstand, onder de voorwaarde, dat alle ε_t onafhankelijk normaal verdeeld zijn met verwachting 0 en variantie σ^2 . Daartoe berekende hij de waarschijnlijkheid, dat zich op $t+1$ een "upcross" zou voordoen, d.w.z. dat geldt

$$(8) \quad x_t = \sum_0^t \alpha^{t-j} \varepsilon_j = \sum_{-1}^t \alpha^{t-j} \varepsilon_j < 0$$

$$(9) \quad x_{t+1} = \sum_0^{t+1} \alpha^{t+1-j} \varepsilon_j = \sum_{-1}^t \alpha^{t-j} \varepsilon_{j+1} > 0,$$

waarbij we ter vereenvoudiging van de formules $\varepsilon_{-1} \equiv 0$ gesteld hebben. Deze waarschijnlijkheid is gelijk aan de massa van

$$(10) \quad dF = K \exp \left[-\frac{1}{2\sigma^2} \sum_0^{t+1} \varepsilon_j^2 \right] d\varepsilon_0 \dots d\varepsilon_{t+1}$$

(K is een constante), die ligt tussen de hypervlakken $\sum_{-1}^t \alpha^{t-j} \varepsilon_j = 0$ en $\sum_{-1}^t \alpha^{t-j} \varepsilon_{j+1} = 0$; immers, de projectie van dF op het $\varepsilon_1, \dots, \varepsilon_{t+1}$ -hypervlak is alzijdig symmetrisch. Deze massa is evenredig met de hoek θ_t tussen eerstgenoemde hypervlakken; hierbij geldt

$$(11) \quad \cos \theta_t = \frac{\sum_{-1}^{t+1} \alpha^{t-j} \alpha^{t-j-1}}{\sum_{-1}^t (\alpha^{t-j})^2} = \frac{1 - \alpha^{2(t+1)}}{1 - \alpha^{2(t+2)}} \alpha$$

Hieruit volgt:

$$(12) \quad \lim_{t \rightarrow \infty} \cos \theta_t = \alpha$$

²⁾ The advanced theory of statistics, II, p. 381 e.v.

Op dezelfde wijze vindt men, als men de veronderstelling $t_0 = 0$ laat varen, voor iedere constante t :

$$(13) \quad \lim_{t_0 \rightarrow -\infty} \cos \theta_t = \alpha.$$

De asymptotische waarschijnlijkheid van een upcross is dus

$$(14) \quad p = \frac{\arccos \alpha}{2\pi}$$

Hieruit concludeert Kendall, dat de verwachting van de afstand tussen twee opeenvolgende "upcrosses" gelijk is aan $2\pi/\arccos \alpha$. Wij zullen trachten dit te interpreteren.

Voeg aan de reeks x_t een tweede reeks y_t toe met de eigenschap, dat als x_t voor $t = t_1$ een "upcross" geeft, y_{t_1} de waarde 1 heeft, terwijl $y_t = 0$ is voor alle waarden van t , die deze eigenschap niet hebben. Dan hebben alle y_t een simultane verdeling met de eigenschap $E y_t = p$. Wanneer we nu een tijdreeks x_t van de lengte n beschouwen, dan krijgen we

$$(15) \quad \lim_{n \rightarrow \infty} E \frac{1}{n} \sum_{t=1}^n y_t = p.$$

Wanneer we het aantal "upcrosses" in de reeks x_t noemen: d_n , dan geldt dus

$$(16) \quad \lim_{n \rightarrow \infty} E n^{-1} d_n = p.$$

Noem verder de gemiddelde afstand tussen opvolgende "upcrosses" \bar{l}_n , dan geldt

$$(17) \quad \lim_{n \rightarrow \infty} \bar{l}_n = \lim_{n \rightarrow \infty} n \cdot d_n^{-1}.$$

Nu is d_n^{-1} een positief-convexe functie van d_n voor $d_n > 0$, d.w.z. dat iedere koorde van deze kromme boven de kromme ligt.³⁾ In dat geval zegt de ongelijkheid van Jensen:

$$(19) \quad E d_n^{-1} \geq (E d_n)^{-1} = \frac{1}{np},$$

zodat $\lim_{n \rightarrow \infty} E \bar{l}_n \geq p^{-1}$ is.

Kendall gaat zelfs verder, en concludeert dat

$$(20) \quad \lim_{n \rightarrow \infty} E \bar{l}_n = \lim_{n \rightarrow \infty} E n d_n^{-1} = p^{-1}$$

is. Dat deze limiet $= p^{-1}$ is, is echter niet bewezen. In het geval, dat de opeenvolgende y_t onderling onafhankelijk zijn, is de overeen-

³⁾ Cf. D. van Dantzig, *Capita Selecta der Waarschijnlijkheidsrekening*. Caput II: Momenten en Ongelijkheden; pp. 93-96.

komstige stelling wel te bewijzen.⁴⁾ Dit is hier echter zeker niet het geval, daar de mogelijkheden $y_t = 1$ en $y_{t+1} = 1$ elkaar zelfs uitsluiten. Hoewel de stelling van Kendall dus niet bewezen is, is echter wel te verwachten, dat zij juist zal blijken te zijn.

6. Ter vergelijking zullen we de verdeling voor de tijdreeks met genoemde stochastische differentievergelijking bestuderen onder de voorwaarde, dat ε_t onafhankelijk homogeen verdeeld is in het interval $(-\frac{1}{2}, \frac{1}{2})$ ⁵⁾, d.w.z. dat geldt:

$$(21) \quad \begin{aligned} P[\varepsilon_t \leq s] &= 0 \quad \text{voor } s \leq -\frac{1}{2} \\ &= \frac{1}{2} + s \quad \text{voor } -\frac{1}{2} \leq s \leq \frac{1}{2} \\ &= 1 \quad \text{voor } s \geq \frac{1}{2}. \end{aligned}$$

We nemen verder aan, dat $|\alpha|$ zo klein is, dat α^2 en hogere machten van α verwaarloosd kunnen worden. In dat geval geldt bij benadering⁶⁾:

$$(22) \quad x_t = \varepsilon_t + \alpha \varepsilon_{t-1}$$

Dan beschouwen we een tweetal termen in de reeks, bijv. x_0 en x_1 , waarvan we veronderstellen, dat de eerste positief, de tweede negatief is. In dat geval begint met x_1 , een iteratie ("run") van termen met negatief teken. Wij zullen nu waarschijnlijkheid berekenen, dat deze iteratie minstens de lengte k heeft, d.w.z. dat geldt

$$(23) \quad x_0 > 0, x_1 < 0, x_2 < 0, \dots, x_k < 0.$$

⁴⁾ Prof. van Dantzig deelt mede, dat zijn uitspraak (Sur la méthode des fonctions génératrices, Coll. de Lyon, 1949, 26-46) "La distribution approximativement" (l.c. p. 42) niet juist is. Met de l.c. gevolgde methode kan bewezen worden, dat voor $k = 2$, $p_1 = p$, $p_2 = q = 1 - p$ deze verwachting gelijk is aan

$$p^{-1} \left[1 + \frac{p^{n+2} - q^{n+2}}{p - q} \right],$$

dus voor n tot p^{-1} nadert mits $pq \neq 0$ is.

⁵⁾ Voor de verdeling van de afstand tussen "upcrosses" maakt het uiteraard geen verschil, of ε_t homogeen verdeeld is tussen $-a$ en $+a$, dan wel tussen $-\frac{1}{2}$ en $+\frac{1}{2}$.

⁶⁾ Alle verdere vergelijkingen in deze paragraaf gelden behoudens machten van de tweede en hogere graad in α .

In symbolen luidt deze waarschijnlijkheid

$$(24) \quad P \left[x_2 < 0, \dots, x_k < 0 \mid x_0 > 0, x_1 < 0 \right]$$

Volgens het vermenigvuldigingsaxioma is zij gelijk aan

$$(25) \quad \frac{P[x_0 > 0, x_1 < 0, \dots, x_k < 0]}{P[x_0 > 0, x_1 < 0]}$$

De noemer hiervan is gelijk aan

$$(26) \quad \begin{aligned} & P \left[\varepsilon_1 + \alpha \varepsilon_0 < 0, \varepsilon_0 + \alpha \varepsilon_{-1} > 0 \right] = \\ & = \int_{-\frac{1}{2}}^{\frac{1}{2}} d\varepsilon_{-1} \int_{-\alpha \varepsilon_{-1}}^{\frac{1}{2}} d\varepsilon_0 \left(\frac{1}{2} - \alpha \varepsilon_0 \right) = \\ & \frac{1}{4} - \frac{\alpha}{8} . \end{aligned}$$

Vervolgens de teller van (20). Stel $k = 2$. Dan geldt:

$$(27) \quad \begin{aligned} & P \left[\varepsilon_2 + \alpha \varepsilon_1 < 0, \varepsilon_1 + \alpha \varepsilon_0 < 0, \varepsilon_0 + \alpha \varepsilon_{-1} > 0 \right] = \\ & = \int_{-\frac{1}{2}}^{\frac{1}{2}} d\varepsilon_{-1} \int_{-\alpha \varepsilon_{-1}}^{\frac{1}{2}} d\varepsilon_0 \int_{-\frac{1}{2}}^{-\alpha \varepsilon_0} d\varepsilon_1 \left(\frac{1}{2} - \alpha \varepsilon_0 \right) = \\ & = \int_{-\frac{1}{2}}^{\frac{1}{2}} d\varepsilon_{-1} \int_{-\alpha \varepsilon_{-1}}^{\frac{1}{2}} d\varepsilon_0 \left(\frac{1}{4} + \frac{\alpha}{8} - \frac{\alpha}{2} \varepsilon_0 \right) \end{aligned}$$

Voor willekeurige positieve gehele k geldt algemeen:

$$(28) \quad \begin{aligned} & P \left[\varepsilon_k + \alpha \varepsilon_{k-1} < 0, \dots, \varepsilon_0 + \alpha \varepsilon_{-1} > 0 \right] = \\ & = \int_{-\frac{1}{2}}^{\frac{1}{2}} d\varepsilon_{-1} \int_{-\alpha \varepsilon_{-1}}^{\frac{1}{2}} d\varepsilon_0 (a_k + b_k - c_k \alpha \varepsilon_0) = \\ & = \frac{1}{2} a_k + \left(\frac{1}{2} b_k - \frac{1}{8} c_k \right) \alpha . \end{aligned}$$

Hierbij is

$$(29) \quad a_k = \frac{1}{2} a_{k-1}$$

$$(30) \quad b_k = \frac{1}{2} b_{k-1} + \frac{1}{8} c_{k-1}$$

$$(31) \quad c_k = a_{k-1}$$

Hieruit en uit de waarden $a_1 = \frac{1}{2}$, $b_1 = 0$, $c_1 = 1$ (zie (21)) volgt:

$$(32) \quad a_k = \left(\frac{1}{2} \right)^k$$

$$(33) \quad b_k = \left(\frac{1}{2} \right)^{k+1} (k-1)$$

$$(34) \quad c_k = \left(\frac{1}{2} \right)^{k-1}$$

Dus is de waarschijnlijkheid, dat de iteratie minstens de lengte k heeft

$$(35) \quad P[\underline{l}' \geq k] = \frac{(\frac{1}{2})^{k+1} + (\frac{1}{2})^{k+2} (k-2)\alpha}{\frac{1}{4} - \frac{\alpha}{8}} =$$

$$= \frac{2 + (k-1)\alpha}{2^k}$$

De complementsfunctie (gedefinieerd als 1 minus de verdelingsfunctie) van \underline{l}' wordt gegeven door: (36) $P[\underline{l}' > k] = \frac{1 + \frac{1}{2} k \alpha}{2^k}$;
de frequentie door:

$$(37) \quad P[\underline{l}' = k] = \frac{1 + (\frac{1}{2} k - 1)\alpha}{2^k}$$

De verwachting van \underline{l}' is

$$(38) \quad \mathcal{E} \underline{l}' = \sum_{k=1}^{\infty} P[\underline{l}' \geq k] = 2 + \alpha.$$

Voor de variantie vindt men

$$(39) \quad \text{var}(\underline{l}') = 2 + 3\alpha.$$

7. In de vorige paragraaf beschouwden we de afstanden van "downcross" tot "upcross"; noem de eerste dezer afstanden \underline{l}'_1 , de tweede \underline{l}'_2 enz. Uit symmetrie-overwegingen volgt, dat de afstanden van "upcross" tot "downcross" $\underline{l}''_1, \underline{l}''_2, \dots$ dezelfde verdeling hebben. Dan geldt voor een "gehele periode" van "upcross" tot "upcross" of van "downcross" tot "downcross" $\underline{l}_i = \underline{l}'_i + \underline{l}''_i$:

$$(40) \quad \mathcal{E} \underline{l}_i = 4 + 2\alpha + \mathcal{O}(\alpha^2)$$

Wanneer de reeks nu bestaat uit m afstanden van "upcross" tot "upcross", dan geldt voor de gemiddelde afstand, dat zijn voorwaardelijke verwachting bij gegeven m gegeven wordt door:

$$(41) \quad \mathcal{E}(\underline{l} | m) = \mathcal{E} \underline{l}_i.$$

We geven de voorwaardelijke verwachting bij gegeven n van de gemiddelde afstand van "upcross" tot "upcross" in geval van normaal verdeelde $\underline{\varepsilon}_t$ aan door $\lambda_N^{(n)}$ en haar limiet voor $n \rightarrow \infty$ door λ_N . Verder geven we de voorwaardelijke verwachting bij gegeven m van genoemde afstand in geval van homogeen verdeelde $\underline{\varepsilon}_t$ aan door λ_H . Dan geldt, indien Kendall's formule juist is:

$$(42) \quad \lambda_N = 4 + \frac{8}{\pi} \alpha + \mathcal{O}(\alpha^2)$$

$$(43) \quad \lambda_H = 4 + 2\alpha + \mathcal{O}(\alpha^2).$$

Hieruit volgt (althans voor voldoende kleine waarden van $|\alpha|$) bij eenzelfde waarde van α , dat λ_N meer afwijkt van de verwachting 4 in geval van onafhankelijk verdeelde x_t dan λ_H . ⁷⁾ Men moet in het oog houden, dat de verwachtingen λ_N en λ_H onder verschillende voorwaarden berekend zijn. Het lijkt ons echter plausibel, dat de betekenis van deze verschillende voorwaarden asymptotisch voor $n \rightarrow \infty$ verdwijnt.

⁷⁾ Is Kendall's formule onjuist, dan moet het gelijkheidsteken in de vergelijking (42) vervangen worden door een groter-teken, zodat deze conclusie gehandhaafd blijft.